

自由語検索のための高速文字列検索方式

2N-8

川口 久光 加藤 寛次 藤澤 浩道 畠山 敦 藤縄 雅章

(株)日立製作所 中央研究所

1. はじめに

文書情報検索システムでは、文書の検索に際して、ユーザが目的の文書について思い付いた自由な言葉（自由語）で検索できる検索システムが求められている。このためには、同じ意味を表わす同義語および表記の異なる異表記語を考慮した検索を行う必要がある。

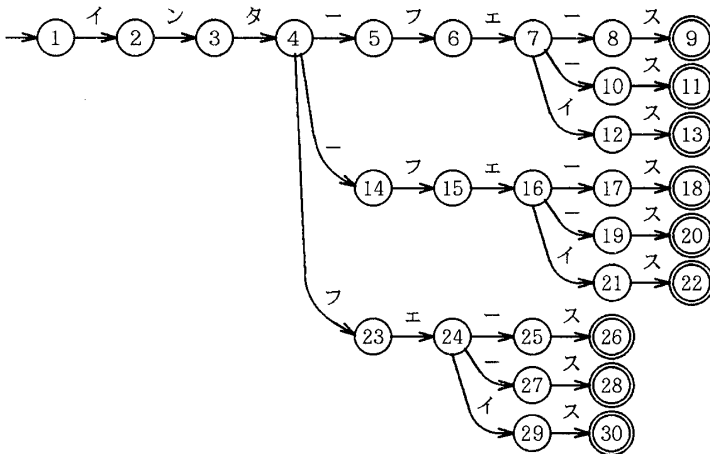
これを実現する方式として、ユーザが投入した検索語から同義語・異表記語を自動的に生成し、これら数百の検索語を一括して照合する自由語検索方式を提案している。[1] この自由語検索において、必要とされる複数検索語の一括検索方法として有限オートマトンを用いたハードウェアを開発してきた。[2]

本稿では、高度な検索機能を実現しようとするオートマトンの状態数が多くなるという従来方式の欠点を改善したコンカレントステート型オートマトン方式を提案するとともに、そのハードウェア方式について言及する。

2. コンカレントステート型オートマトン方式

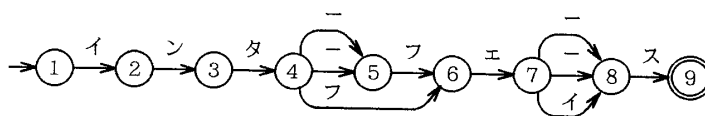
従来は、活性化しているオートマトンの状態を示すトークンを1つだけ照合処理に使用していた。このため、状態遷移のフェイル処理を考慮してオートマトンを作成しなければならないという制約があった。この制約のため木状に広がるオートマトンしか記述出来ず、異表記語を照合するオートマトンを作成しようとする状態数が非常に多くなってしまふという問題があった（第1図参照）。

この課題を解決するために、複数のトークンを照合処理に使用することによりフェイル処理を不要としたコンカレントステート型オートマトン（CSA）方式を提案する。本方式ではフェイル処理が不要なため、オートマトンの状態遷移を網状にまとめて記述することができる。その結果、異表記用のオートマトンも極めて少ない状態数で実現することが可能となった（第2図参照）。



第1図 異表記検索用オートマトン（従来方式）

イン $\begin{bmatrix} \text{タ} \\ \text{フ} \end{bmatrix}$ $\begin{bmatrix} \text{エ} \\ \text{ス} \end{bmatrix}$ ス : 複合語表現形式



テキスト:

イ	ン	タ	フ	エ	イ	ス
---	---	---	---	---	---	---

 トークン: 1 → 2 → 3 → 4 → 6 → 7 → 8 → 9
 の移動: 1 → 2

第2図 異表記検索用オートマトン（CSA方式）

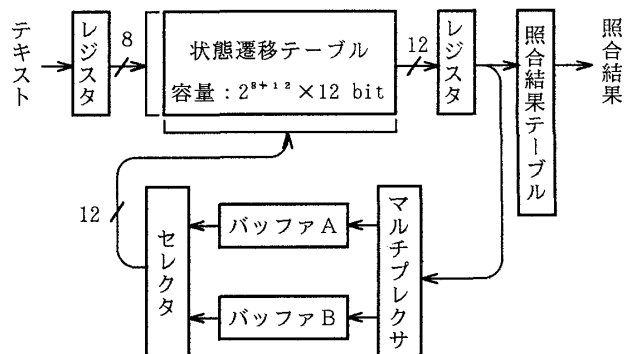
A Fast String Searching Algorithm for Free Term Search

Hisamitsu KAWAGUCHI Kanji KATO Hiromichi FUJISAWA Atsushi HATAKEYAMA Masaaki FUJINAWA
 HITACHI, Ltd.

3. CSA型サーチエンジンの構成

CSA方式に基づく検索ハードウェア（サーチエンジン）は、オートマトンを格納する状態遷移テーブル、トークンを格納する2面バッファ、および照合結果識別子を格納する照合結果テーブルから構成する（第3図参照）。第3図では、4096状態のオートマトンを格納することが可能である。

1回の状態遷移は、テキストから1文字とバッファAからトークンが存在する現在の状態番号を読み込み、状態遷移テーブルを参照してからトークンが移動すべき次の状態番号を読み出しバッファBに格納することにより実現している。バッファAとバッファBを切り替えながら、次々と文字コードを読み込み同様の処理を繰り返すことによって文字列の照合を行う。



第3図 CSA方式サーチエンジンの構成

4. 高度な検索機能

CSA方式では、従来困難であった下記①～④の高機能な検索を少ない状態数のオートマトンで実現できる。

① 異表記検索：

インタフェース → インタ~~ニ~~フエース
 インタ~~ニ~~フエ~~イ~~ス
 ……
 ……

② don't care文字指定検索：

漢字??方式 → 漢字~~認~~識方式

③ 近傍条件指定検索：

文字[3C]方式 → 文字~~認~~識方式
 文字~~読~~取り方式

④ 1文字誤り許容検索：

読取り方式 → 読~~み~~取り方式
 読~~取~~方式

5. おわりに

本方式によれば、フェイル処理が不要なためオートマトンの状態遷移を網状にまとめて記述でき、従来に比べ少ない状態数で検索オートマトンを作成することが可能となった。このことにより、異表記検索、don't care文字指定検索、近傍条件指定検索、1文字誤り許容検索などの高度な検索機能を有する高速サーチエンジンを実現できるようになった。

参考文献

- [1] 島山他：“知的ファイリングシステムの開発（その2）；自由語検索における異表記、異表現解消法”，第33回情処全大 4Y-9（1986）。
- [2] 川口他：“高速検索文字列検索方式”，昭62信学全大54（1987）。